

Approximate Dynamic Programming: a Q -Function Approach

Paul Beuchat, Angelos Georghiou and John Lygeros¹

Abstract—In this paper we study both the value function and Q -function formulation of the Linear Programming (LP) approach to ADP. The approach selects from a restricted function space to fit an approximate solution to the true optimal Value function and Q -function. Working in the discrete-time, continuous-space setting, we extend and prove guarantees for the fitting error and online performance of the policy, providing tighter bounds. We provide also a condition that allows the Q -function approach to be more efficiently formulated in many practical cases.

Index Terms—Approximate Dynamic Programming (ADP), Bellman Inequality, Iterated Bellman Inequality, Q functions, Distributed Control

I. INTRODUCTION

Stochastic optimal control problems are solved exactly by the methodology of Dynamic Programming (DP). In 1952 Bellman proposed a solution method for discrete time problems involving general dynamics and cost function, [1]. The solution of the Bellman Equation is the optimal cost-to-go function which characterizes the behaviour of the optimal control policy. Although a powerful result, computing a solution to the Bellman equation introduces many challenges. Namely that the solution lives in an infinite dimensional function space, and multi-variate expectations can be neither analytically expressed nor tractably computed. The continuous state, input, and disturbance spaces could be discretized to leverage the extensive literature on solving the Bellman equation for discrete spaces, eg. [3], [18]. This paper focuses on continuous space problems of sufficiently high dimension for which discretization method are not feasible.

Function approximation for Approximate Dynamic Programming (ADP) is a method which restricts the space of functions considered so that an approximate, sub-optimal, solution of the Bellman equation can be tractably found. The Linear Programming (LP) approach to ADP is a method that approximates the cost-to-go function, [20]. In this paper we focus on a variant that approximates instead the Q -function. The Q -function, first introduced in [25], is similar to the cost-to-go function, but it has the property that the optimal control policy can be written in terms of the optimal Q -function without involving any of the terms that describe the model. This property of the Q -function makes the formulation interesting for tackling decentralized control problems, first suggested in [8] for discrete space problems. A drawback of

the Q -function formulation is that it doubles the size of the LP that needs to be solved.

A key challenge of the LP approach to ADP is choosing a metric that indicates which function from the restricted function space yields a high quality approximation. This challenge exists equally for the Value function and the Q -function formulation. In [9], the authors presented a variant of the LP approach that allowed them to give theoretical guarantees on the quality of the approximation for discrete space problems. They provided two types of guarantees: (i) a performance bound on the online performance of the control policy, (ii) a fitting bound on how closely the approximate Value function is to the true optimal. It is non-trivial to extend the performance bound to continuous space because the continuous space transition kernel cannot be directly inverted. One of the fitting bounds was tightened slightly in [24] by using an iterated version of the Bellman Inequality. However, this doesn't provide additional insight for how to compute a higher quality approximation.

Many other solution methods approximate the stochastic optimal control problem, other than the Dynamic Programming reformulation, see [17] for a survey. Two well known examples are Receding Horizon Model Predictive Control (MPC) [19], [6] and Linear Decision Rules [2], [11]. In this paper we provide a link between these two approximations and the LP approach to ADP.

The contributions contained in this paper are:

- We introduce the iterated Q -function formulation of the LP approach to ADP and provide a condition for when it can be solved more efficiently by eliminating half the decision variables and constraints.
- We prove the online performance bound for continuous space. This bound justifies that the LP approach is sensible for continuous space applications.
- We propose an iterated version of the greedy policy and bound the sub-optimality of its online performance. The iterated greedy policy indicates a link between ADP and MPC solution methods.
- We use the iterated Bellman inequality to tighten the fitting error bound that is based on Lyapunov functions.

The structure of this paper is as follows. Section II presents the Dynamic Programming formulation considered. Section III introduces the Q -function approximation methods to be studied. Section IV provides the theoretical performance guarantees for both the Value function and Q -function formulations. Section V provides the condition for when the Q -function formulation can be simplified. Section VI uses numerical examples to demonstrate the theory and bounds, and Section VII concludes.

This research was partially funded by the European Commission under the project Local4Global.

¹All Authors are with the Automatic Control Laboratory, Department of Information Technology and Electrical Engineering, ETH Zürich, Switzerland {beuchatp, angelosg, jlygeros}@control.ee.ethz.ch

Notation: \mathbb{R}_+ is the space of non-negative scalars; \mathbb{Z}_+ is the space of positive integers; I_n is the $n \times n$ identity matrix; $(\cdot)^\top$ is the matrix transpose; given $f: \mathcal{X} \rightarrow \mathbb{R}$, the infinity norm is $\|f\|_\infty = \sup_{x \in \mathcal{X}} |f(x)|$, and the weighted 1-norm is $\|f\|_{1,c} = \int_{\mathcal{X}} |f(x)| c(dx)$.

The term intractable is used throughout the paper. We loosely define intractable to mean that the computational burden of any existing solution method prohibits finding a solution in reasonable time.

II. DYNAMIC PROGRAMMING (DP) FORMULATION

This section introduces the problem formulation and states the DP as the solution to the Bellman equation for both Value functions and Q -functions. We then formulate a LP whose solution is equivalent to the Bellman equation. The LP motivates the ADP approach presented in Section III.

A. Bellman Equation and Operator

We consider infinite horizon, stochastic optimal control problems with a discounted cost objective. The system is described by discrete dynamics over continuous state and action spaces. The state of the system at time t is denoted by $x_t \in \mathcal{X} \subseteq \mathbb{R}^{n_x}$. The system state is influenced by control decisions $u_t \in \mathcal{U} \subseteq \mathbb{R}^{n_u}$, and stochastic disturbances $\xi_t \in \Xi \subseteq \mathbb{R}^{n_\xi}$. In this setting, the state evolves according to the function $g: \mathcal{X} \times \mathcal{U} \times \Xi \rightarrow \mathcal{X}$ as, $x_{t+1} = g(x_t, u_t, \xi_t)$. At time t , the system incurs the stage cost $\gamma^t l(x_t, u_t)$, where $\gamma \in [0, 1)$ is the discount factor and the objective is to minimize the infinite sum of the stage costs.

The optimal Value function, $V^*: \mathcal{X} \rightarrow \mathbb{R}$, characterizes the solution of the stochastic optimal control problem. It represents the cost-to-go from any state of the system if the optimal control policy is played. The optimal Value function is the solution of the Bellman equation [1],

$$V^*(x) = \min_{u \in \mathcal{U}} \underbrace{l(x, u) + \gamma \mathbb{E}[V^*(g(x, u, \xi))]}_{(\mathcal{T}V^*)(x)}, \quad (1)$$

which holds for all $x \in \mathcal{X}$. The operator \mathcal{T} is the well known Bellman operator, and the \mathcal{T}_u operator is equivalent but without the minimization over u .

The optimal Q -function, $Q^*: (\mathcal{X} \times \mathcal{U}) \rightarrow \mathbb{R}$, characterizes the solution of the stochastic optimal control problem. It represents the cost of making decision u now and then playing optimally from the next time step forward. The optimal Q -function is the solution of the following variant of the Bellman equation:

$$\begin{aligned} Q^*(x, u) &= (\mathcal{T}_u V^*)(x, u), \\ &= l(x, u) + \gamma \mathbb{E} \left[\underbrace{\min_{v \in \mathcal{U}} Q^*(g(x, u, \xi), v)}_{(FQ^*)(x, u)} \right], \end{aligned} \quad (2)$$

which holds for all $x \in \mathcal{X}$, and $u \in \mathcal{U}$. The F -operator is the equivalent of \mathcal{T} , but instead for Q -functions.

The optimal policy is generated from either V^* or Q^* ,

$$\pi^*(x) = \arg \min_{u \in \mathcal{U}} l(x, u) + \gamma \mathbb{E}[V^*(g(x, u, \xi))], \quad (3a)$$

$$= \arg \min_{u \in \mathcal{U}} Q^*(x, u). \quad (3b)$$

This is commonly referred to as the *Greedy Policy*. Note that evaluating (3a) uses the dynamics, stage cost, and expectation with respect to ξ , whereas (3b) involves only Q^* .

B. LP Reformulation of DP

In this sub-section, we formulate an LP whose optimal solution is the same Q^* that solves equation (2). We develop the formulation in terms of Q -functions, the same line of reasoning holds for Value function, see [24].

The *Bellman Inequality* is a well known relaxation of the Bellman equation for Value functions. Equivalently, equation (2) can also be relaxed to an inequality,

$$Q(x, u) \leq FQ(x, u), \quad \forall x \in \mathcal{X}, u \in \mathcal{U}, \quad (4)$$

which will be referred to as the *F-operator inequality*. As operator F operator is monotone, and satisfies value iteration convergence, any Q satisfying (4) will be a point-wise under-estimator of Q^* .

Let $\mathcal{F}(\mathcal{X} \times \mathcal{U})$ denote the function space of real-valued measurable functions on $(\mathcal{X} \times \mathcal{U})$ with finite weighted ∞ -norm. Under the same assumptions as [12, §6.3], it follows that the solution of the following LP,

$$\begin{aligned} \max_Q \quad & \int_{\mathcal{X} \times \mathcal{U}} Q(x, u) c(dx, du) \\ \text{s.t.} \quad & Q \in \mathcal{F}(\mathcal{X} \times \mathcal{U}), \\ & Q(x, u) \leq FQ(x, u), \quad \forall x \in \mathcal{X}, u \in \mathcal{U}, \end{aligned} \quad (5)$$

satisfies (2), when $c(\cdot, \cdot)$ is a finite measure on $(\mathcal{X} \times \mathcal{U})$ that assigns a positive mass to all open subsets of $(\mathcal{X} \times \mathcal{U})$. The equivalence between (2) and (5) requires that $\mathcal{F}(\mathcal{X} \times \mathcal{U})$ is used as the function space over which the decision variable, Q , is optimized, see [12, §6.3]. Intuitively speaking, the reason is that the space $\mathcal{F}(\mathcal{X} \times \mathcal{U})$ is rich enough to satisfy $Q \leq FQ$ with equality, point-wise for all $x \in \mathcal{X}$ and $u \in \mathcal{U}$.

The constraint in (5) is not linear in Q due to the minimization inside the expectation. The following proposition linearizes the constraint by introducing an additional decision variable. Letting $\mathcal{F}(\mathcal{X})$ denote the function space of real-valued measurable functions on \mathcal{X} with finite weighted ∞ -norm, we state the following equivalent reformulation of the F -operator inequality.

Proposition 2.1: For an arbitrary $Q: (\mathcal{X} \times \mathcal{U}) \rightarrow \mathbb{R}$ the following are equivalent:

- (i) $Q(x, u) \leq FQ(x, u)$
- (ii) $\exists V \in \mathcal{F}(\mathcal{X}) : Q(x, u) \leq \mathcal{T}_u V(x), V(x) \leq Q(x, u)$

where the inequalities hold for all $x \in \mathcal{X}$ and $u \in \mathcal{U}$.

Proof: See [8, Theorem 2]. ■

Note that (i) \Leftrightarrow (ii) provided that $V \in \mathcal{F}(\mathcal{X})$. If V is taken to be in some subset of $\mathcal{F}(\mathcal{X})$, then the reformulation is only sufficient, i.e., (ii) \Rightarrow (i).

C. Iterated Bellman Inequality

The feasible region of (5) can be increased by using an iterated F -operator inequality. The iterated formulation has the same optimizer and optimal value as (5). For the iterated Value function formulation, see [24].

Any Q -function satisfying an iterated F -operator inequality, $Q \leq F^M Q$ will be a point-wise under-estimator of Q^* . This follows from the monotone and value iteration convergence property of F . The same reasoning as [12, §6.3] also establishes that the solution of the following LP,

$$\begin{aligned} \max_Q \quad & \int_{\mathcal{X} \times \mathcal{U}} Q(x, u) c(d(x, u)) \\ \text{s.t.} \quad & Q \in \mathcal{F}(\mathcal{X} \times \mathcal{U}), \\ & Q(x, u) \leq F^M Q(x, u), \quad \forall x \in \mathcal{X}, u \in \mathcal{U}, \end{aligned} \quad (6)$$

satisfies (2), when $c(\cdot, \cdot)$ is chosen in the same manner. For $M = 1$, (5) and (6) are equivalent.

The effect of using $M > 1$ is to increase the feasible region of under-estimators. However, the feasible region of (6) does not strictly increase as a function of M , see [24, §3.4]. The same feasible set relationship holds when (6) is approximated by a using restricted function space. See Figure 2 for a schematic representation

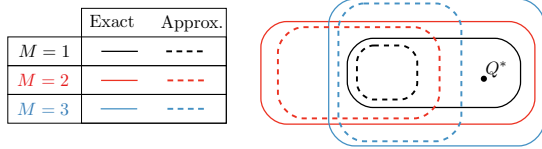


Fig. 1: Showing the feasible region of (6) for increasing number of iterations, M , of the iterated F -operator inequality. The feasible set for (6) is represented by the solid-lines and contains Q^* for all values of M . The dotted-lines depict a case where (6) is restricted to a function space that does not contain the optimal.

The iterated F -operator inequality constraint is not linear in Q due to multiple nested minimizations and expectations. The constraint can be reformulated as M separate F -operator inequalities by introducing $M-1$ additional decision variables as per the following proposition. Using Proposition 2.1 to linearize each F -operator inequality, (6) becomes an LP.

Proposition 2.2: For an arbitrary $Q : (\mathcal{X} \times \mathcal{U}) \rightarrow \mathbb{R}$ the following are equivalent:

- (i) $Q(x, u) \leq F^M Q(x, u)$
- (ii) $\exists Q_1, \dots, Q_{M-1} \in \mathcal{F}(\mathcal{X} \times \mathcal{U})$ such that:

$$\begin{aligned} Q(x, u) &\leq F Q_1(x, u), \\ Q_{j-1}(x, u) &\leq F Q_j(x, u), \quad j = 2, \dots, M-1, \\ Q_{M-1}(x, u) &\leq F Q(x, u), \end{aligned}$$

where the inequalities hold for all $x \in \mathcal{X}$ and $u \in \mathcal{U}$.

Proof: Follows from [24, §3.4]. ■

Similar to the reformulation of the F -operator, if for any j , Q_j is taken to be in some subset of $\mathcal{F}(\mathcal{X} \times \mathcal{U})$, then the reformulation is only sufficient, i.e., (ii) \Rightarrow (i).

D. Sources of Intractability

Solving (6) for Q^* , and implementing (3b), is in general intractable. The difficulties can be categorized as:

- (D1) $\mathcal{F}(\mathcal{X})$ and $\mathcal{F}(\mathcal{X} \times \mathcal{U})$ are infinite dimensional spaces;
- (D2) Problem (6) has infinite constraints;
- (D3) The objective of (6) is a multidimensional integral;
- (D4) The multidimensional integral over ξ in the F -operator;
- (D5) For arbitrary $Q^* \in \mathcal{F}(\mathcal{X} \times \mathcal{U})$, the greedy policy (3b) may be intractable;

Difficulties (D1-D5) apply equally to the Value function formulation and represent *curses of dimensionality*, see [16].

III. APPROXIMATE DYNAMIC PROGRAMMING (ADP)

In this section we restrict the function spaces $\mathcal{F}(\mathcal{X})$ and $\mathcal{F}(\mathcal{X} \times \mathcal{U})$ to simultaneously overcome (D1-D5), and hence make (6) computationally tractable. The solution of the approximate LP is then used to implement the greedy policy.

A. The Approximate LP

As suggested in [20], we restrict the admissible Value functions and Q -functions to be expressed as a linear combination of basis functions. In particular, given basis functions $\hat{V}^{(i)} : \mathcal{X} \rightarrow \mathbb{R}$ and $\hat{Q}^{(i)} : (\mathcal{X} \times \mathcal{U}) \rightarrow \mathbb{R}$, we parameterize the restricted function spaces as,

$$\begin{aligned} \hat{\mathcal{F}}(\mathcal{X}) &= \left\{ \sum_{i=1}^K \alpha_i \hat{V}^{(i)}(x) \mid \alpha_i \in \mathbb{R}, \right\}, \\ \hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U}) &= \left\{ \sum_{i=1}^K \beta_i \hat{Q}^{(i)}(x, u) \mid \beta_i \in \mathbb{R} \right\}. \end{aligned} \quad (7)$$

Hence an element of either set is fully specified by a choice of α_i 's or β_i 's. An approximate solution of (6) is obtained by using these restricted function spaces in the following *approximate iterated LP*:

$$\begin{aligned} \max_{\hat{Q}} \quad & \int_{\mathcal{X} \times \mathcal{U}} \hat{Q}(x, u) c(d(x, u)) \\ \text{s.t.} \quad & \hat{Q} \in \hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U}), \\ & \hat{Q}(x, u) \leq F^M \hat{Q}(x, u), \quad \forall x \in \mathcal{X}, u \in \mathcal{U}, \end{aligned} \quad (8)$$

where the only change from (6) was to replace $\mathcal{F}(\mathcal{X} \times \mathcal{U})$ by $\hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$. The optimization variables are now the β_i 's in the definition of $\hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$. To make the constraint in (8) linear, Proposition 2.1 and 2.2 are used with all the additional Value functions and Q -functions restricted to $\hat{\mathcal{F}}(\mathcal{X})$ and $\hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$ respectively. It is possible to use a different basis function set for each additional Value function and Q -function.

Given a solution of (8), denoted \hat{Q}^* , a natural choice for the online policy is, to replace Q^* in equation (3), i.e.,

$$\hat{\pi}(x) = \arg \min_{u \in \mathcal{U}} \hat{Q}^*(x, u) \quad (9)$$

called an *approximate greedy policy*.

In general, Q^* is not an element of the restricted function space. Thus \hat{Q}^* will not be the solution of the Bellman equation (2). The following lemma provides the intuition that \hat{Q}^* is the closest under-estimator to Q^* , relative to the choice of $c(\cdot, \cdot)$. See III-C below for more discussion.

Lemma 3.1: An approximate Q -function, \hat{Q} , solves (8), if and only if it solves a minimization problem with the same decision variables, same constraints, and the objective replaced by $\|Q^* - \hat{Q}\|_{1,c(x,u)}$.

Proof: See [9, Lemma 1] \blacksquare

Difficulty (D1) has been overcome for problem (8) as $\hat{\mathcal{F}}(\mathcal{X})$ and $\hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$ are parameterized by a finite dimensional decision variable.

B. Options for overcoming (D2-D5)

There are a number of choices of $\hat{\mathcal{F}}(\mathcal{X})$ and $\hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$ that efficiently address (D2-D5). The possible choices depend on the functional form of the stage cost and dynamics, the description of the state and input sets, the distribution of the exogenous disturbance, and the method used to overcome (D2). Table I summarises a range examples found in the literature.

TABLE I: Examples of overcoming (D2-D5)

Ref.	Problem instance studied:	Overcome (D2) by:
[24], [22]	Polynomial problems	S-procedure
[10], [14], [23]	Various	Sampling
[13]	Stochastic reachability	Sampling
[15]	Perimeter surveillance	Exact Reformulation

C. Choice of relevance weighting

In the exact LP, (6), the specific choice of $c(\cdot, \cdot)$ does not affect the optimal solution, as long as it satisfies a mild restriction. This is no longer the case in (8) where the choice of $c(\cdot, \cdot)$ plays a central role in determining the quality of \hat{Q}^* . Lemma 3.1 suggests that in regions where $c(\cdot, \cdot)$ is higher, \hat{Q}^* is expected to give an under-estimate that is closer to the true optimal, Q^* . Thus, $c(\cdot, \cdot)$ is commonly referred to as the *relevance weighting*.

A good approximation of the optimal Q -function could be described as one for which the online performance of the approximate greedy policy is near optimal. Although Lemma 3.1 shows that \hat{Q}^* is the closest approximate Q -function for a given set of basis functions, it says nothing about the sub-optimality of playing policy (9). In Section IV we show that the online performance of playing (9) can be bounded by how well \hat{Q}^* approximates Q^* . In fact, the weighting parameter in the bound gives a hint for how to choose the relevance weighting needed in the approximate LP.

It is possible to alleviate the inherent compromise of only having a single choice for $c(\cdot, \cdot)$. In [4] the authors suggest solving (8) for multiple choices of $c(\cdot, \cdot)$, and using the point-wise maximum from the family of approximations in the greedy policy. They argue that improved online performance can be achieved with their approach.

D. Choice of M

Any $M > 1$ has an increased feasible region relative to $M = 1$. Increasing M increases the size of decision variables and constraints in (8), and hence increases the solve time. It is advisable to choose M as large as possible such that the approximate LP can be solved in the time frame available.

E. Improved Approximate Policy

Applying the F -operator to \hat{Q}^* will give an improved approximation of Q^* and possibly yield an improved approximate greedy policy. The value iteration convergence property of F means that $(F\hat{Q})(x, u)$ is a better approximate of Q^* than \hat{Q} , in an ∞ -norm sense. Thus, given an approximate Q -function, \hat{Q} , the *iterated policy*,

$$\hat{\pi}(x) = \arg \min_{u \in \mathcal{U}} F^D \hat{Q}(x, u) \quad (10)$$

improves on (9), with $D \geq 1$ as the number of iterations.

The complication is that even when $\hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$ was chosen to make (9) tractable, this improved policy is not tractable because of the nested expectations and minimizations arising from the $F^D \hat{Q}$ term. By similar arguments, an improved choice of policy in terms of an approximate Value function, \hat{V} , would be

$$\hat{\pi}(x) = \arg \min_{u \in \mathcal{U}} l(x, u) + \mathbb{E} \left[\left(\mathcal{T}^D \hat{V} \right) (f(x, u, \xi)) \right], \quad (11)$$

which also involves nested expectations and minimizations.

Writing out the iterations of the F or \mathcal{T} operator, it can be seen that this improved policy is exactly the generic form of a D -stage stochastic programming problem [21, section 3.1]. A popular approximate solution method for such stochastic programs is Model Predictive Control (MPC). Solving (10) or (11) with an MPC approach would be equivalent to a finite horizon MPC formulation, with a time horizon of D steps, and using \hat{Q} or \hat{V} as the terminal cost.

In Section IV, we give a bound on the sub-optimality of the online performance achieved by (10) or (11). This indicates that a tighter performance bound can be achieved through the improved approximate policy and in Section VI we use a numerical example to demonstrate the potential.

IV. PERFORMANCE BOUNDS FOR ADP

In this section, we present the infinite space performance bounds for both Value functions and Q -functions. The three types of bounds presented we present were first presented for discrete-space Value functions in [9]. Our contribution in this section is two-fold: (i) we show that the bounds hold in continuous spaces for both Value functions and Q -functions, (ii) we present an improved version for the online performance bound and an improved version of the Lyapunov-based bound. The online performance bound justifies the improved approximate policy suggested in section III-E,

A. Online Performance Bound

We present first a bound on the online performance of playing the improved approximate policy (10) or (11). These bounds do not depend on what method was used to compute the approximate Value function, they only require that \hat{V} is a point-wise under-estimator of V^* .

Before presenting the bounds, we introduce two measures: the expected state-by-action frequency, μ defined on $(\mathcal{X} \times \mathcal{U})$, and its marginal on the state space, $\bar{\mu}$ defined on \mathcal{X} , called

the expected state frequency. For any Borel set $\Gamma \in \mathcal{B}(\mathcal{X} \times \mathcal{U})$ and $B \in \mathcal{B}(\mathcal{X})$ the measures are defined as:

$$\begin{aligned} \mu(\Gamma) &:= \sum_{t=0}^{\infty} \gamma^t P_{\nu}^{\pi}[(x_t, u_t) \in \Gamma] \\ \tilde{\mu}(B) &:= \mu(B \times \mathcal{U}) = \sum_{t=0}^{\infty} \gamma^t P_{\nu}^{\pi}[x_t \in B] \end{aligned} \quad (12)$$

where the overloaded notation $P_{\nu}^{\pi}[(\cdot)]$ represents the probability given that the initial states are distributed according to ν and the system evolves autonomously under the fixed policy π . Intuitively speaking, $P_{\nu}^{\pi}[(\cdot)]$ can be seen as simulating the autonomous system from samples in the initial distribution, and keeping count of how many times the set of interest, Γ or $(B \times \mathcal{U})$, is visited. See [12, 6.3.6] for further details.

Lemma 4.1: $(1 - \gamma)\tilde{\mu}$ is a probability distribution.

Proof: Evaluating the measure over \mathcal{X} ,

$$\begin{aligned} (1 - \gamma)\tilde{\mu}(\mathcal{X}) &= (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \underbrace{P_{\nu}^{\pi}[x_t \in \mathcal{X}]}_{=1} \\ &= (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t = 1 \end{aligned}$$

To simplify the presentation of the theorems, let us define some notation. Given any policy $\pi : \mathcal{X} \rightarrow \mathcal{U}$, let $V_{\pi} : \mathcal{X} \rightarrow \mathbb{R}$ denote the online performance, i.e., $V_{\pi}(x)$ is the cost-to-go from state x when policy π is played. By definition V_{π} is a point-wise over-estimator of V^* . Additionally, given an function $Q : (\mathcal{X} \times \mathcal{U}) \rightarrow \mathbb{R}$, define the following: $Q|_{\pi}(x) := Q(x, \pi(x))$.

Theorem 4.2: Let $\hat{Q} : (\mathcal{X} \times \mathcal{U}) \rightarrow \mathbb{R}$ be such that $\hat{Q}(x, u) \leq Q^*(x, u)$ for all $x \in \mathcal{X}$, $u \in \mathcal{U}$, and let $\hat{\pi} : \mathcal{X} \rightarrow \mathcal{U}$ be a D -iterated policy defined as per (10). Then the sub-optimality of the online performance is bounded as,

$$\|V_{\hat{\pi}} - V^*\|_{1, \nu} \leq \frac{1}{1 - \gamma} \left\| Q^*|_{\hat{\pi}} - (F^D \hat{Q})|_{\hat{\pi}} \right\|_{1, (1 - \gamma)\tilde{\mu}}$$

Proof: See appendix A. ■

Theorem 4.3: Let $\hat{V} : \mathcal{X} \rightarrow \mathbb{R}$ be such that $\hat{V}(x) \leq V^*(x)$ for all $x \in \mathcal{X}$, and let $\hat{\pi} : \mathcal{X} \rightarrow \mathcal{U}$ be a D -iterated policy defined as per (11). Then the sub-optimality of the online performance is bounded as,

$$\|V_{\hat{\pi}} - V^*\|_{1, \nu} \leq \frac{1}{1 - \gamma} \left\| V^* - (\mathcal{T}^D \hat{V}) \right\|_{1, (1 - \gamma)\tilde{\mu}}$$

Proof: Follows as a minor adaptation of proof given for Theorem 4.2 ■

When $D = 0$, Theorems 4.2 and 4.3 agree with the discrete space versions, [8, Theorem 1] and [9, Theorem 1] respectively. Figure 2 visualizes the quantities involved in the online performance bound.

The following insights apply to Theorem 4.2 and 4.3:

- They provide the reassurance for continuous space problems that by playing policy (10) or (11) using an under-estimator function, the sub-optimality of the

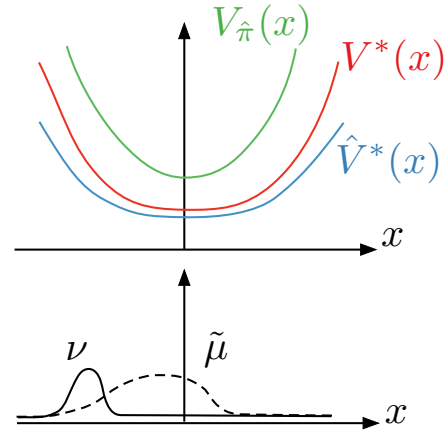


Fig. 2: The upper axes show that the online performance of playing policy $\hat{\pi}$ is a point-wise over-estimator of V^* . By assumption of Theorem 4.3, \hat{V} is a point-wise under-estimator of V^* and hence $\mathcal{T}\hat{V}$ is also. The lower axis highlights that 1-norm weightings in Theorem 4.3, ν and $\tilde{\mu}$, can be very different.

online performance is bounded by how closely \hat{Q} or \hat{V} fits Q^* or V^* respectively.

- They motivate the potential benefit of considering a D -iterated policy based on an under-estimator function. Although F and \mathcal{T} are not contractive w.r.t. the weighted 1-norm, it is expected that the RHS gets smaller as D increases, and hence the online sub-optimality is more tightly bounded.
- Lemma 3.1 showed that the solution of the approximate LP (8) minimizes effectively the right-hand-side of the bound with $c(\cdot, \cdot)$ as the 1-norm weighting. Hence the bound suggests that the expected state frequency, $\tilde{\mu}$ is a natural choice for the relevance weighting. This, however, is circular, because $\tilde{\mu}$ depends on the solution of the approximate LP.

In the discrete space setting, an analytic expression is derived for the expected state frequency. As all the quantities in discrete space are represented by vectors and matrices, the infinite sum in (12) can be written as an inverse of the transition kernel matrix and greatly simplifies that proof.

In the next sub-section we present two theorems that further bound the right-hand-side of Theorems 4.2 and 4.3.

B. Infinty-norm Fitting Bound

We present now the first result that bounds the fitting of \hat{Q}^* or \hat{V}^* by how close Q^* or V^* is to the span of the basis functions. The bounds consider all functions from the restricted function space. This type of bound was first presented in [9] for discrete space, and the main contribution of [24] was to improve this bound by using the iterated Value function formulation. In [4] the bound was proven for the Q -function formulation.

Theorem 4.4: Given Q^* is the solution of (2), and \hat{Q}^* is the solution of (8) for a given choice $\hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$ and $c(\cdot, \cdot)$,

then the following bound holds,

$$\|Q^* - \hat{Q}^*\|_{1,c(x,u)} \leq \frac{2}{1-\gamma^M} \min_{\hat{Q} \in \hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})} \|Q^* - \hat{Q}\|_\infty$$

Proof: See Appendix B ■

Theorem 4.5: Given V^* is the solution of (1), and \hat{V}^* is the solution of the approximate iterated LP for a given choice $\hat{\mathcal{F}}(\mathcal{X})$ and $c(\cdot)$, then the following bound holds,

$$\|V^* - \hat{V}^*\|_{1,c(x)} \leq \frac{2}{1-\gamma^M} \min_{\hat{V} \in \hat{\mathcal{F}}(\mathcal{X})} \|V^* - \hat{V}\|_\infty$$

Proof: See [24, §4.3] ■

Comparing the left-hand-side of Theorem 4.5 to the right-hand-side of Theorem 4.3, we see that if $c(\cdot)$ is chosen to be the expected state frequency, then the online performance is also bounded by Theorem 4.5. For Theorems 4.4 and 4.2 to be combined in a similar way, the relevance weighting $c(\cdot, \cdot)$ should satisfy,

$$\|Q^*|_{\hat{\pi}} - (F^D \hat{Q})|_{\hat{\pi}}\|_{1,\hat{\mu}} = \|Q^* - \hat{Q}^*\|_{1,c(x,u)}.$$

In both cases, this is again a circular requirement because the choice of the relevance weighting affects the solution of the approximate LP, which affects the approximate policy, which affects the expected state frequency, which in turn affects the choice of relevance weighting.

The following insights apply to Theorem 4.4 and 4.5:

- The terms $\|Q^* - \hat{Q}\|_\infty$ and $\|V^* - \hat{V}\|_\infty$ can be overwhelming large, even for the minimum taken over the span of the basis function set. Practically, it may not be possible to choose M large enough when solving the approximate LP to overcome this.
- The right-hand-side of the bounds hold for any choice of the relevance weighting. Thus, the bounds do not provide any intuition for how to choose the relevance weighting, $c(\cdot, \cdot)$ and $c(\cdot)$, to achieve the tightest fit.

In the next sub-section we present two bounds that depend on the relevance weighting.

C. Lyapunov-based Fitting Bound

The fitting bound can be improved for solutions of the iterated approximate LP by using the Lyapunov functions defined in this sub-section. This Lyapunov-based fitting bound follows the lines of [9] and tightens the bound based on the iterated Bellman inequality. In [4, §IV-D] it was shown that the infinity-norm bounds apply also to the case of using a point-wise maximum of Q -functions. The same holds true for the Lyapunov-based bounds presented in this sub-section.

In order to give the Lyapunov function definitions, we introduce the following two operators. For any Value function, $V : \mathcal{X} \rightarrow \mathbb{R}$, define operator H_V as,

$$(H_V V)(x) = \max_{u \in \mathcal{U}} \mathbb{E}[V(f(x, u, \xi))],$$

and for any $Q : (\mathcal{X} \times \mathcal{U}) \rightarrow \mathbb{R}$, define the operator H_Q as,

$$(H_Q Q)(x, u) = \max_{v \in \mathcal{U}} \mathbb{E}[Q(f(x, u, \xi), v)].$$

Given that the system is in state x , the function $(H_V V)(x)$ represents the worst case expected value of the next state. For Q -functions, given further that action u will be applied, the function $(H_Q Q)(x, u)$ represents the worst case expected value two times steps into the future given that the system evolves stochastically to state $f(x, u, \xi)$. It is readily shown that both the H_V and H_Q operators are monotone.

For any Value function or Q -function, the following scalars are defined:

$$\beta_V = \max_{x \in \mathcal{X}} \frac{\gamma (H_V V)(x)}{V(x)},$$

$$\beta_Q = \max_{(x,u) \in (\mathcal{X} \times \mathcal{U})} \frac{\gamma (H_Q Q)(x, u)}{Q(x, u)},$$

which both represent the maximum ratio of: the worse case expected value at a future time step, with the value in the current state(-by-input).

Thus the notion of *Lyapunov functions* that will be used in the theorems are defined as follows.

Definition 4.6: A Value function $V : \mathcal{X} \rightarrow \mathbb{R}_{++}$ is called a Lyapunov Value function if $\beta_V < 1$.

Definition 4.7: A Q -function $Q : (\mathcal{X} \times \mathcal{U}) \rightarrow \mathbb{R}_{++}$ is called a Lyapunov Q -function if $\beta_Q < 1$.

For any strictly positive function, $V : \mathcal{X} \rightarrow \mathbb{R}_{++}$, let $1/V$ denote the map $x \mapsto 1/V(x)$, and similarly for a strictly positive Q -function. Now we can state the improved bounds.

Theorem 4.8: Given V^* is the solution of (1), and \hat{V}^* is the solution of the approximate iterated LP. Then, for any function $\hat{V}^+(x) \in \hat{\mathcal{F}}(\mathcal{X})$ that is a Lyapunov Value function as per definition 4.6, the following bound holds,

$$\|V^* - \hat{V}^*\|_{1,c(x)} \leq \frac{2 \|\hat{V}^+\|_{1,c(x)}}{1 - \beta_{\hat{V}^+}^M} \min_{\hat{V} \in \hat{\mathcal{F}}} \|V^* - \hat{V}\|_{\infty, 1/\hat{V}^+}$$

Proof: See appendix C. ■

Theorem 4.9: Given Q^* is the solution of (2), and \hat{Q}^* is the solution of (8). Then, for any function $\hat{Q}^+(x, u) \in \hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$ that is a Lyapunov Q -function as per definition 4.7, the following bound holds,

$$\|Q^* - \hat{Q}^*\|_{1,c(x,u)} \leq \frac{2 \|\hat{Q}^+\|_{1,c(x,u)}}{1 - \beta_{\hat{Q}^+}^M} \min_{\hat{Q} \in \hat{\mathcal{F}}} \|Q^* - \hat{Q}\|_{\infty, 1/\hat{Q}^+}$$

Proof: The proof follows the same steps as the proof of Theorem 4.8, adapted accordingly. ■

The following insights apply to Theorem 4.8 and 4.9:

- As the infinity norm of the right-hand-side of the bounds is weighted by the inverse of the Lyapunov function, it will generally be tighter than the equivalent term in Theorems 4.4 and 4.5. To see this, consider that in regions where V^* or Q^* are large, the Lyapunov function should also be large and hence reduce the worst case error in those regions.

- The relevance weighting now appears on the right-hand side of the bound. This indicates that another appropriate choice of relevance weighting is that which gives the tightest bound. However, finding the combination of a relevance weighting and Lyapunov function that yields the tightest bound is, in general, a difficult problem.

See [9, §5] for some discussion on the choice of Lyapunov functions for discrete space problems.

D. Road Map to Bounds

To assist the reader, the bounds presented above, and those referred to in [8], [24], and [9] are summarised in Table II. The four new bounds presented in this section represent our contribution to the Road Map.

TABLE II: Road Map to Bounds

Online performance bound		
Space	Value functions	Q-functions
Discrete	[9, Theorem 1]	[8, Theorem 1]
Continuous	Theorem 4.3	Theorem 4.2

Infinity-norm bound		
Space	Value functions	Q-functions
Discrete	[9, Theorem 2]	Follows from [9, Theorem 2]
Continuous	[24, §4.2]	[4, Theorem 4.1]

Lyapunov-based bound		
Space	Value functions	Q-functions
Discrete	[9, Theorem 1]	Follows from [9, Theorem 1]
Continuous	Theorem 4.3	Theorem 4.2

V. AN EFFICIENT Q-FUNCTION FORMULATION

In this section, we analyze a case where the Q-function formulation can be made more efficient by eliminating constraints and decision variables. The improved formulation has the same number of infinite constraints and decision variables as the Value function formulation. We present two classes of problems to highlight that the efficient formulation is applicable in many practical situations.

A. Condition for equivalence

The iterated approximate Q-function formulation from Section III-A allows each decision variable to be taken from a different restricted function space. Letting $\hat{\mathcal{F}}_0(\mathcal{X} \times \mathcal{U}) \subseteq \mathcal{F}(\mathcal{X} \times \mathcal{U})$ denote an alternative basis function set, and applying Proposition 2.1 and 2.2 to problem (8), leads to,

$$\begin{aligned}
& \max_{\hat{Q}_0, \hat{V}_j} \int_{\mathcal{X} \times \mathcal{U}} \hat{Q}_0(x, u) c(d(x, u)) \\
& \text{s.t. } \hat{Q}_0, \hat{Q}_M \in \hat{\mathcal{F}}_0(\mathcal{X} \times \mathcal{U}) \\
& \quad \hat{Q}_j \in \hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U}), \quad j = 1, \dots, M-1, \\
& \quad \hat{V}_j \in \hat{\mathcal{F}}(\mathcal{X}), \quad j = 0, \dots, M-1, \\
& \quad \hat{Q}_j(x, u) \leq \mathcal{T}_u \hat{V}_j(x, u), \quad j = 0, \dots, M-1, \quad (13a) \\
& \quad \hat{V}_j(x) \leq \hat{Q}_{j+1}(x, u), \quad j = 0, \dots, M-1, \quad (13b) \\
& \quad \hat{Q}_0 = \hat{Q}_M, \quad (13c)
\end{aligned}$$

where the inequality constraints hold for all $x \in \mathcal{X}$ and $u \in \mathcal{U}$. Now consider the following formulation with $M-1$ fewer Q-functions and M fewer infinite constraints:

$$\begin{aligned}
& \max_{\hat{Q}_0, \hat{V}_j} \int_{\mathcal{X} \times \mathcal{U}} \hat{Q}_0(x, u) c(d(x, u)) \\
& \text{s.t. } \hat{Q}_0 \in \hat{\mathcal{F}}_0(\mathcal{X} \times \mathcal{U}) \\
& \quad \hat{V}_j \in \hat{\mathcal{F}}(\mathcal{X}), \quad j = 0, \dots, M-1, \\
& \quad \hat{Q}_0(x, u) \leq \mathcal{T} \hat{V}_0(x, u), \quad (14a) \\
& \quad \hat{V}_{j-1}(x) \leq \mathcal{T} \hat{V}_j(x, u), \quad j = 1, \dots, M-1, \quad (14b) \\
& \quad \hat{V}_{M-1}(x) \leq \hat{Q}_0(x, u), \quad (14c)
\end{aligned}$$

where the inequality constraints hold for all $x \in \mathcal{X}$ and $u \in \mathcal{U}$. In Lemma 5.1 below, we provide a condition for when (13) and (14) are equivalent.

Lemma 5.1: If sets $\hat{\mathcal{F}}(\mathcal{X})$ and $\hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$ satisfy that for all $\hat{V} \in \hat{\mathcal{F}}(\mathcal{X})$ there exists a $\hat{Q} \in \hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$ such that

$$\hat{Q}(x, u) = \mathcal{T}_u \hat{V}(x, u), \quad \forall x \in \mathcal{X}, u \in \mathcal{U},$$

then the approximate LP (13) is equivalent to (14), in the sense that both problems have the same optimal value and there is a mapping between feasible and optimal solutions in both problems.

Proof: See Appendix D ■

The condition of this lemma gives an indication of when it is unnecessary to formulate problem (13), problem (14) being preferred because of the fewer decision variables and constraints. Problem (14) can be seen as fitting an approximate Q-function to the constraints of the iterated Value function formulation. Without the assumption of Lemma 5.1, the constraints of (14) do not in general imply that the constraints of (13) can be satisfied and hence do not imply that $\hat{Q}_0 \leq F^M \hat{Q}_0$.

In the remainder of this section we provide two examples of classes of problems where the condition is satisfied.

B. Input constrained, Linear-Quadratic

In the case of linear dynamics, quadratic cost function, and control actions constrained to lie in a polytopic feasible set, then the optimal Value function and Q-function is known to be piece-wise quadratic [5, Theorem 7.7]. Hence quadratic basis functions are a reasonable choice for both $\hat{\mathcal{F}}(\mathcal{X})$ and $\hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$, see [24, §6]. We will show that, in this setting, the condition of Lemma 5.1 is satisfied and hence formulation (14) is preferred over (13).

The quadratic basis functions used for the restricted function spaces are defined as,

$$\begin{aligned}
\hat{\mathcal{F}}(\mathcal{X}) &= \left\{ \hat{V}(x) \mid \begin{array}{l} V(x) = x^\top P x + p^\top x + s \\ P \in \mathbb{S}^{n_x}, p \in \mathbb{R}^{n_x}, s \in \mathbb{R} \end{array} \right\} \\
\hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U}) &= \left\{ \hat{Q}(x, u) \text{ such that: } \begin{array}{l} Q(x, u) = \begin{bmatrix} x \\ u \end{bmatrix}^\top P_Q \begin{bmatrix} x \\ u \end{bmatrix} + p_Q^\top \begin{bmatrix} x \\ u \end{bmatrix} + s_Q \\ P_Q \in \mathbb{S}^{n_x+n_u}, p_Q \in \mathbb{R}^{n_x+n_u}, s_Q \in \mathbb{R} \end{array} \right\}
\end{aligned}$$

where the α_i 's and β_i 's from (7) are the coefficients of the monomials. In this setting, for any quadratic Value function, the term

$$\mathcal{T}_u \hat{V}(x, u) = l(x, u) + \mathbb{E}[\hat{V}(f(x, u, \xi))],$$

will be quadratic in $[x^\top, u^\top]^\top$ and require knowledge of the first and second moments of the exogenous disturbance. As $\hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$ is taken to be the space of all quadratic functions in $[x^\top, u^\top]^\top$, the condition of Lemma 5.1 is satisfied.

C. Structured \mathcal{Q} functions for decentralized control

Now we consider an example motivated by decentralized control. In a decentralized control problem, the input and state vectors are split up so that the computation of each input can only depend on a certain portion of the state vector. By contrast, exact dynamic programming assumes that full state information is available at each time step. Decentralized control is realized if evaluation of the approximate greedy policy, (9), has the separable structure required. This dictates a structure for the approximate \mathcal{Q} -function.

Before describing the required structure of a \mathcal{Q} -function, let us first introduce the notation of a decentralized control problem involving N agents. Let $u = [u_1^\top, \dots, u_N^\top]^\top$ be a partition of the control action vector that each agent needs to decide, and let x_1, \dots, x_N denote the portion of the state vector available to the respective agent for making its decision. Then the decentralized policy should consist of N separate *local* policies, each depending only on the relevant portion of the state vector, i.e.,

$$u = \pi_{\text{Decent}}(x) = \begin{bmatrix} \pi_1(x_1) \\ \vdots \\ \pi_N(x_N) \end{bmatrix}.$$

The greedy policy (9) requires the evaluation of a constrained optimisation problem. Consequently, for the greedy policy to be separable, both the objective and constraint set must have the required separable structure. The constraint enforced by the greedy policy is $u \in \mathcal{U}$, thus the input constraint set is assumed to be separable, i.e.,

$$\mathcal{U} = \mathcal{U}_1 \times \dots \times \mathcal{U}_N.$$

The \mathcal{Q} -function is the objective of the greedy policy (9). Hence, the objective is separable if the \mathcal{Q} -function is a sum of per-agent \mathcal{Q} -functions that only depend on the control decisions to be made by that agent and the portion of the state vector available. Let \mathcal{S} denote the set of functions of $[x^\top, u^\top]^\top$ possessing the following separable structure,

$$\mathcal{S} = \left\{ \hat{Q}(\cdot, \cdot) \left| \hat{Q}(x, u) = \sum_{i=1}^N \hat{Q}_i(u_i, x_i) + q(x) \right. \right\} \quad (15)$$

where $q(x)$ can be any function of the full state vector. The term $q(x)$ is allowed because it does not affect the decision made by evaluating the greedy policy. This separable \mathcal{Q} -function structure was first suggested in [8].

This ADP approach to decentralized control is not formulated to solve the true decentralized control problem. It is

instead formulated to fit an approximate \mathcal{Q} -function to the centralized optimal in such a way that the resulting policy is decentralized.

Given the solution from either (13) or (14), \hat{Q}_0^* is the \mathcal{Q} -function to be used in the approximate greedy policy. Thus, it is only necessary to enforce the constraint (15) on \mathcal{Q}_0 . The remaining \mathcal{Q} -functions and Value functions in the iterated approximate LP should not have the decentralized structure enforced as it adds unnecessary constraints. Lemma 5.1 allows for a different restricted function space for \mathcal{Q}_0 , hence, the lemma can be applied to the decentralized control formulation. For a decentralized, input constrained, linear-quadratic problem with quadratic basis function sets, we enforce structure $\hat{Q}_0 \in \mathcal{S}$, and allow $\{\hat{V}_j, \hat{Q}_j\}_{j=1}^{M-1}$ to be dense quadratics. This satisfies the condition of Lemma 5.1, and hence the additional $M-1$ \mathcal{Q} -functions can be eliminated without changing the solution.

The Value function formulation can also be used to approximate a solution to the decentralized control problem. A separable Value function based greedy policy requires the assumption that $l(x, u) \in \mathcal{S}$. When solving the approximate LP, $\hat{\mathcal{F}}(\mathcal{X})$ should be further restricted so that the approximate Value function satisfies $\mathbb{E}[\hat{V}(f(x, u, \xi))] \in \mathcal{S}$.

VI. NUMERICAL RESULTS

In this section we present two numerical examples to highlight various aspects of the theory from Sections II to V. The first example exemplifies the potential of the improved approximate policy, and the second numerical example demonstrates that using a structured \mathcal{Q} -function for decentralized control can achieve near centralized performance. Both examples highlight the difficulty in selecting the state relevance weighting.

A. 1-dimensional example

We use the example in this sub-section to illustrate that the improved approximate policy, (11), does shift the policy closer to the optimal, and we highlight that, although the iterated value function gives an improved lower bound, it has worse online performance. The 1-dimensional example we use is taken directly from [24], it has size $n_x = n_u = n_\xi = 1$. The dynamics, costs, and constraints are given by,

$$\begin{aligned} x_{t+1} &= x_t - 0.5u_t + \xi, & l(x, u) &= x^2 + 0.1u^2, \\ \gamma &= 0.95, & |u| &\leq 1 \end{aligned}$$

with the exogenous disturbance and initial condition distributed as $\xi_t \sim \mathcal{N}(0, 0.1)$ and $x_0 \sim \mathcal{N}(0, 10) = \nu$ respectively. The benefit of using a one-dimensional example is that the true optimal Value Function and optimal policy (V^* and π^*) can be computed by using a discretisation method and plotted for gaining visual insight. Additionally, the online performance of each policy can be readily computed via Monte-Carlo simulation.

We use the space of univariate quadratics as $\hat{\mathcal{F}}(\mathcal{X})$, and will chosen the state-relevance weighting to be the initial state distribution, i.e., $c(dx) = \nu(dx)$. Relative to the optimal, we will compare two approximate value functions,

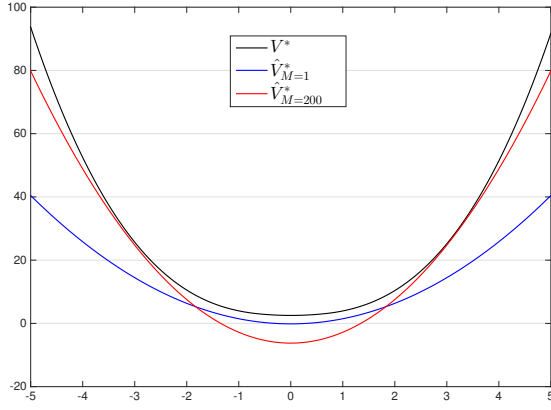


Fig. 3: Showing the value functions for the 1-dimensional example. The approximate value functions $\hat{V}_{M=1}^*$ (blue) and $\hat{V}_{M=200}^*$ (red) are compared to the optimal V^* (black).

solved via the iterated approximate LP with $M = 1$ and $M = 200$, denoted $\hat{V}_{M=1}^*$ and $\hat{V}_{M=200}^*$ respectively. These choices are inline with [24].

The approximate greedy policy will be compared with its improved version, (11). Using $\hat{V}_{M=1}^*$ and $\hat{V}_{M=200}^*$, we denote the approximate greedy policies as $\hat{\pi}_{M=1}$ and $\hat{\pi}_{M=200}$ respectively. We will use $\hat{V}_{M=200}^*$ to play the improved policy with 3 iterations, denoted $\hat{\pi}_{M=200}^{D=3}$. This policy is implemented by approximating the solution of (11) with as 3 step MPC problem, using $\mathbb{E}[\xi] = 0$ as the predicted disturbance.

Figure 3 shows the two approximate value functions compared to the optimal, while figure 4 shows the three approximate policies compared to the optimal. Table III compares the optimal cost-to-go with the two lower-bounds made from the approximate value functions, and with the online performance of the approximate policies. All the values are with expectation taken over the initial distribution, $x_0 \sim \nu$.

TABLE III: Results for the 1-dimensional example

Description	Value
Online performance of $\hat{\pi}_{M=200}$	37.84
Online performance of $\hat{\pi}_{M=1}$	37.81
Online performance of $\hat{\pi}_{M=200}^{D=3}$	37.81
Optimal value function V^*	37.80
Lower bound computed from $\hat{V}_{M=200}^*$	28.20
Lower bound computed from $\hat{V}_{M=1}^*$	16.09

It can be seen in figure 3 that $\hat{V}_{M=1}^*$ (blue) and $\hat{V}_{M=200}^*$ (red) are both point-wise lower-bounds of V^* (black). Table III shows that $\hat{V}_{M=200}^*$ gives an improved lower-bound of the optimal. Although, we see in figure 3 that $\hat{V}_{M=1}^*$ gives a better point-wise lower-bound in the region near $x = 0$, it is clear that $\hat{V}_{M=200}^*$ gives a significantly better lower-bound in the regions further from $x = 0$. Thus when the expectation is taken over x_0 , $\hat{V}_{M=200}^*$ more tightly lower-bounds V^* .

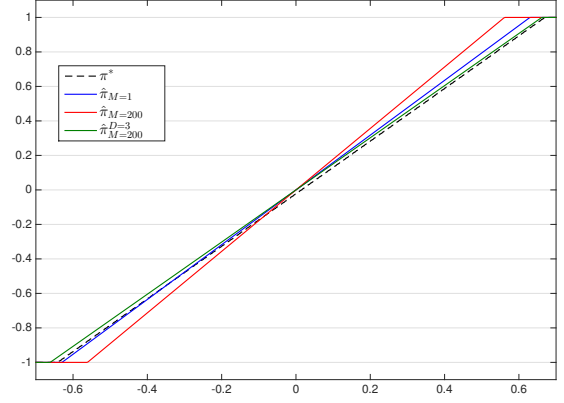


Fig. 4: Showing the policies for the 1-dimensional example. Outside of the \mathcal{X} region shown all the policies saturate at ± 1 . The approximate greedy policies $\hat{\pi}_{M=1}$ (blue) and $\hat{\pi}_{M=200}$ (red), and the improved approximate policy $\hat{\pi}_{M=200}^{D=3}$ (green) are compared to the optimal π^* (black dashed).

Despite $\hat{V}_{M=1}^*$ giving a worse lower-bound of the optimal value function, it can be seen in figure 4 that the policy it generates $\hat{\pi}_{M=1}$ (blue) is a significantly better approximate of the optimal policy π^* (black dashed) as compared to with $\hat{\pi}_{M=200}$ (red). The reason is that the greedy policy is determined by the gradient of the value function. Therefore, in regions where the gradient of a \hat{V} closely approximates V^* , then the greedy policies will generate closely matching control actions. In figure 3 it is clear that in the region near the origin $\hat{V}_{M=1}^*$ matches the gradient of V^* much better than $\hat{V}_{M=200}^*$. Due to the input constraints of this problem, outside of that region all value functions that rise steeply enough play the same because the input saturates at ± 1 .

The difficulty in choosing the state-relevance weighting is highlighted by the fact that the approximate value function with $M = 200$ gives a better lower-bound of V^* but has a worse policy and online performance. For the $M = 200$ approximate LP, it would be possible to choose a $c(x)$ different from ν that yields the same value function as $\hat{V}_{M=1}^*$. Thus there is an inherent discrepancy between choosing a $c(x)$ that maximizes the lower-bound of V^* , useful for assessing sub-optimality, and choosing a $c(x)$ that achieves the best online performance, the objective of the stochastic optimal control problem.

The improved approximate policy $\hat{\pi}_{M=200}^{D=3}$, shown in green in figure 4, is a possible method to alleviate the difficulty in choosing the relevance weighting. Figure 4 shows that even with a very modest number of iterations, i.e., $D = 3$, the improved policy significantly improves on its starting point $\hat{\pi}_{M=200}$.

B. Coupled Oscillator Example

In this sub-section we analyze an application of Q -functions for decentralized control. We use a string of coupled oscillators, visualized as a spring-mass-dampener

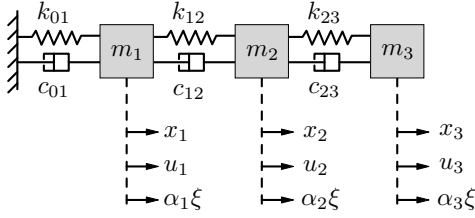


Fig. 5: Schematic showing the coupled oscillator model used for this example to demonstrate using Q -functions for decentralized control. Each mass is considered as a separate constituent system, and needs to make its control decision u_i based only on the measurement of its own state, x_i and \dot{x}_i .

system in figure 5. Each mass is considered as a separate constituent system, and needs to make its control decision u_i based only on the measurement of its own state.

The dynamics for the coupled oscillator shown in Figure 5 are a linear system that can be readily derived by writing the equations of motion for each mass. The state vector has two components for each mass, x_i and \dot{x}_i , which represent position and velocity respectively. Each is controlled by means of the scalar input u_i which represents the driving force applied to the mass. The stochasticity in the system is model by a 1-dimensional uncertainty ξ that represents an exogenous driving force, for example wind, and the scalar factor α_i specifies the influence on each mass. The spring constant and dampening ratio of the mechanical elements connecting mass i to mass j are denoted by k_{ij} and c_{ij} respectively. The fixed wall is represented as $i = 0$.

In order to be able to tractably compute the centralized optimal solution of this larger example, quadratic costs and unconstrained inputs are used, making this an LQR example. The stage cost for each mass is $l_i(x_i, \dot{x}_i, u_i) = x_i^2 + \dot{x}_i^2 + 0.5u_i^2$, with a discount factor of $\gamma = 0.99$ used. The remaining parameters are chosen homogeneously as,

$$m_i = 1, k_{ij} = 3, c_{ij} = 0.05, \alpha_i = 0.1, \forall i, j.$$

The exogenous disturbance and initial condition distributed as $\xi_t \sim \mathcal{N}(0, 1)$ and $x_0 \sim \mathcal{N}(0, I_6) = \nu$ respectively.

The Q -function can produce a decentralized policy to play online if given an appropriate structure, as we describe in section V-C. The coupled oscillator setup here is an LQR example for which the optimal Q -function is quadratic. Due to the coupling in the dynamics, the centralized optimal Q -function will not yield a greedy policy with a separable structure. Hence we use convex quadratic Q -functions as the basis set, $\hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$, further restricted to those having the structure shown in figure 6. Thus, the approximate greedy policy is decentralized, and $Q^* \notin \hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$.

Table IV presents the results of using structured Q -functions for decentralized control of the coupled oscillator system with 3 masses and the parameters given above. We solve (8) with $M = 1$ and $M = 200$, denoted $\hat{Q}_{M=1}^*$ and $\hat{Q}_{M=200}^*$ respectively. Table IV gives the lower-bound of the optimal for each approximate Q -function, and also the online

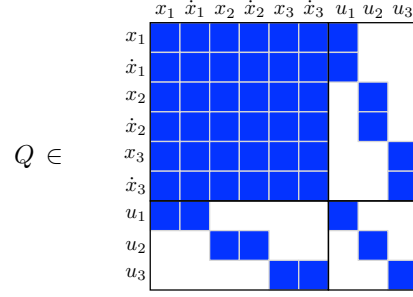


Fig. 6: Showing the quadratic Q -function structure required for the coupled oscillator example so that the greedy policy is decentralized, i.e., control action u_i can be computed from measurements of only x_i and \dot{x}_i . Shaded squares represent order 2 monomial coefficients that can be non-zero.

performance of the approximate greedy policy. The lower-bound is computed as $\mathbb{E}_\nu [\min_u \hat{Q}(x, u)]$, and the online performance is computed via Monte Carlo simulation.

TABLE IV: Results for 3-mass Coupled Oscillator example

Description	Value
Online performance with $M = 200$	138.9
Online performance with $M = 200$ and $D = 2$	133.8
Online performance with $M = 1$	132.2
Online performance with $M = 1$ and $D = 2$	129.4
Optimal value function V^*	128.7
Lower bound computed with $M = 200$	126.6
Lower bound computed with $M = 1$	126.4

The results in table IV show that the decentralized ADP approach using Q -functions, can produce near centralized optimal performance. Specifically for this problem, the online performance using $\hat{Q}_{M=1}^*$ is within 3% of the optimal centralized performance. Interestingly, we see again for this example that although $\hat{Q}_{M=200}^*$ gives a slightly tighter lower bound, the approximate greedy policy generated by $\hat{Q}_{M=1}^*$ yields better online performance, and in this example it is a 4.6% improvement.

Using both $\hat{Q}_{M=1}^*$ and $\hat{Q}_{M=200}^*$ we play also the improved approximate policy with $D = 2$ iterations. The improved approximate policy is no longer decentralized due to the coupled dynamics producing a non-separable structure in the MPC re-formulation of the policy. The results, in table IV, demonstrate that the MPC re-formulation of the improved approximate policy can improve the online performance, by 2 – 3% in this case.

VII. CONCLUSIONS

In this paper we have extended upon existing theoretical results for the Linear Programming approach to Approximate Dynamic Programming. We provided theoretical guarantees that the online performance is bounded when using the approximate greedy policy to make decisions for controlling continuous space systems. The performance guarantees were presented equally for both the value function and the Q -function formulations. We also proposed an improved policy

with tighter theoretical bounds on the online performance, and demonstrated its potential through numerical examples.

We proposed a condition that allows for a significantly more efficient Q -function formulation. As an example of the practical application of this condition, we analyzed a decentralized optimisation methodology based on Q -functions. The key insight is that the approximate greedy policy becomes separable when the Q -function is given the appropriate structure. This methodology was explored through a numerical example of a coupled oscillator where the online performance using decentralized Q -functions achieved within 3% of the optimal centralized performance.

APPENDIX

A. Proof of online performance bound

The proof of Theorem 4.2 uses the machinery presented in [12, section 6.3] for exact DP. Before presenting the proof of Theorem 4.2, we introduce the key objects required and their properties. First, let $K(\cdot|\cdot, \cdot)$ denote the discrete-time transition kernel describing the state evolution under the dynamics and the exogenous and control inputs, i.e., given a set $B \in \mathcal{B}(\mathcal{X})$ then,

$$K(B | x_t, u_t) = \mathbb{P}[f(x_t, u_t, \xi) \in B]$$

represents the probability that state x_{t+1} will be in B given that the system is currently in state x_t and input u_t is played.

The transition kernel allows us to introduce two operators. The first operator, T_π , acts on the space of finite signed measures $\mathcal{M}(\mathcal{X})$. Given a measure $\rho \in \mathcal{M}(\mathcal{X})$, a feasible policy $\pi : \mathcal{X} \rightarrow \mathcal{U}$, and a set $B \in \mathcal{B}(\mathcal{X})$, the operator T_π is defined as,

$$(T_\pi \rho)(B) = \rho(B) - \gamma \int_{x \in \mathcal{X}} K(B | x, \pi(x)) \rho(dx).$$

Thus T_π represents the discounted difference in occupancy measure between two time steps of the stochastic process. The second operator, T_π^* , acts on the space of bounded functions $\mathcal{F}(\mathcal{X})$. Given a function $V \in \mathcal{F}(\mathcal{X})$, and the same feasible policy, the operator T_π^* is defined as,

$$\begin{aligned} (T_\pi^* V)(x) &= V(x) - \gamma \int_{y \in \mathcal{X}} V(y) K(dy | x, \pi(x)) \\ &= V(x) - \gamma \mathbb{E}[V(f(x, \pi(x), \xi))] \end{aligned} \quad (16)$$

Thus T_π^* represents the expected value of discounted difference between two time steps of the stochastic process. Both operators define a continuous linear map and are adjoint,

$$\int_{x \in \mathcal{X}} V(x) (T_\pi \rho)(dx) = \int_{x \in \mathcal{X}} (T_\pi^* V)(x) \rho(dx), \quad (17)$$

see [12, section 6.3]. The online performance bound for discrete space is proven by inverting the transition kernel matrix, see [9, Theorem 1]. The adjoint property of T_π and T_π^* can be seen as a parallel to inverting the transition kernel.

A required identity is that the online performance can be expressed in terms of the stage cost and the frequency

measure defined in Section IV-A. Given a policy, $\pi : \mathcal{X} \rightarrow \mathcal{U}$ and the expected state frequency with respect to that policy, $\tilde{\mu}$, the online performance is expressed as,

$$\begin{aligned} V_\pi(y) &:= \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t l(x_t, \pi(x_t)) \mid x_0 = y \right] \\ &= \int_{\mathcal{X}} l(x, \pi(x)) \tilde{\mu}(dx) \end{aligned} \quad (18)$$

When the left hand side is integrated over the initial state distribution, ν , then $\tilde{\mu}$ is chosen accordingly.

The final identity required relates the initial state distribution to the expected state frequency. Given any $B \in \mathcal{B}(\mathcal{X})$ the following relation holds:

$$\begin{aligned} \nu(B) &= \tilde{\mu}(B) - \gamma \int_{x \in \mathcal{X}} K(B | x, \pi(x)) \tilde{\mu}(dx) \\ &= (T_\pi \tilde{\mu})(B) \end{aligned} \quad (19)$$

This identity stems from [12, eq. (6.3.10)].

We now have all the tools required to prove Theorem 4.2.

Theorem 4.2 Let $\hat{Q} : (\mathcal{X} \times \mathcal{U}) \rightarrow \mathbb{R}$ be such that $\hat{Q}(x, u) \leq Q^*(x, u)$ for all $x \in \mathcal{X}$, $u \in \mathcal{U}$, and let $\hat{\pi} : \mathcal{X} \rightarrow \mathcal{U}$ be a D -iterated policy defined as per (10). Then the sub-optimality of the online performance is bounded as,

$$\|V_{\hat{\pi}} - V^*\|_{1,\nu} \leq \frac{1}{1-\gamma} \left\| Q^*|_{\hat{\pi}} - \left(F^D \hat{Q} \right)|_{\hat{\pi}} \right\|_{1,(1-\gamma)\tilde{\mu}}$$

Proof: By assumption we have for all $k \in \mathbb{N}$,

$$\hat{Q}(x, u) \leq F^k \hat{Q}(x, u) \leq Q^*(x, u) \leq Q_{\hat{\pi}}(x, u), \quad (20)$$

for all $x \in \mathcal{X}$ and $u \in \mathcal{U}$, and hence also for all $u = \hat{\pi}(x) \in \mathcal{U}$.

Recalling the notation $Q|_\pi(x) := Q(x, \pi(x))$, we have,

$$\begin{aligned} &\|V_{\hat{\pi}} - V^*\|_{1,\nu} \\ &= \int_{\mathcal{X}} (V_{\hat{\pi}}(x) - V^*(x)) \nu(dx) \\ &\leq \int_{\mathcal{X}} \left(V_{\hat{\pi}}(x) - \left(F^D \hat{Q} \right)|_{\hat{\pi}}(x) \right) \nu(dx) \\ &= \int_{\mathcal{X}} l(x, \hat{\pi}(x)) \tilde{\mu}(dx) - \int_{\mathcal{X}} \left(F^D \hat{Q} \right)|_{\hat{\pi}}(x) (T_\pi \tilde{\mu})(dx) \\ &= \int_{\mathcal{X}} l(x, \hat{\pi}(x)) \tilde{\mu}(dx) - \int_{\mathcal{X}} \left(T_\pi^* \left(F^D \hat{Q} \right)|_{\hat{\pi}} \right)(x) \tilde{\mu}(dx) \\ &= \int_{\mathcal{X}} \left(F^{D+1} \hat{Q} \right)|_{\hat{\pi}}(x) \tilde{\mu}(dx) - \int_{\mathcal{X}} \left(F^D \hat{Q} \right)|_{\hat{\pi}}(x) \tilde{\mu}(dx) \\ &\leq \int_{\mathcal{X}} Q^*|_{\hat{\pi}}(x) \tilde{\mu}(dx) - \int_{\mathcal{X}} \left(F^D \hat{Q} \right)|_{\hat{\pi}}(x) \tilde{\mu}(dx) \\ &= \frac{1}{1-\gamma} \left\| Q^*|_{\hat{\pi}} - \left(F^D \hat{Q} \right)|_{\hat{\pi}} \right\|_{1,(1-\gamma)\tilde{\mu}} \end{aligned}$$

The first equivalence and first inequality hold by the pointwise ordering of (20). The second equivalence uses (18) for the first term and (19) for the second term. The third equality uses the fact that T_π^* is the adjoint operator of T_π , see (17). The fourth equivalence uses (16) to expand the T_π^* operator,

and then the definition of the F -operator and the chosen policy to construct the first term. The second inequality and final equivalence also hold by the point-wise ordering of (20) and the definition of the 1-norm. The factor $(1 - \gamma)$ was introduced so that the scaling in the 1-norm is a probability distribution. ■

B. Proof of infinity-norm fitting bound for \mathcal{Q} -functions

Theorem 4.4 is the next theorem presented in the paper. The proof requires two additional lemmas that are presented first, and then we present the proof of Theorem 4.4. Lemma 1.1 provides a point-wise bound on how much the M -iterated F -operator inequality is violated for any given \mathcal{Q} function, from the basis or otherwise. This is used in the proof of Lemma 1.2, which shows that given a $\hat{Q} \in \hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$, it can be downshifted by a certain constant amount to satisfy the iterated F -operator inequality. The constant by which it is downshifted relates directly to the constant on the RHS of Theorem 4.4. The proof here is an adaptation to \mathcal{Q} -functions of the proof for Value functions that is given in [24, §4.3].

Lemma 1.1: Let $M \in \mathbb{N}$, and let $Q : (\mathcal{X} \times \mathcal{U}) \rightarrow \mathbb{R}$ be any \mathcal{Q} -function, then violations of the iterated F -operator inequality can be bounded as,

$$(F^M Q)(x, u) \geq Q(x, u) - (1 + \gamma^M) \|Q^* - Q\|_\infty,$$

for all $x \in \mathcal{X}$, $u \in \mathcal{U}$.

Proof: Starting from the terms not involving γ ,

$$\begin{aligned} & Q(x, u) - \|Q^* - Q\|_\infty - (F^M Q)(x, u) \\ & \leq Q^*(x, u) - (F^M Q)(x, u), \quad \forall x \in \mathcal{X}, u \in \mathcal{U} \\ & \leq \| (F^M Q^*) - (F^M Q) \|_\infty \\ & \leq \gamma^M \|Q^* - Q\|_\infty. \end{aligned}$$

The first inequality follows from the definition of the ∞ -norm, and the second inequality comes from $Q^*(x, u) = (FQ^*)(x, u)$ and the ∞ -norm definition. Finally, the third inequality is due to the γ -contractive property of the F -operator. Re-arranging, the result follows. ■

Lemma 1.2: Let $\hat{Q}(x, u) \in \hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$ be an arbitrary element from the basis functions set, and let $\tilde{Q}(x, u)$ be a \mathcal{Q} -function defined as,

$$\tilde{Q}(x, u) = \hat{Q}(x, u) - \underbrace{\frac{1 + \gamma^M}{1 - \gamma^M} \|Q^* - \hat{Q}\|_\infty}_{\text{downwards shift term}}, \quad (21)$$

then $\tilde{Q}(x, u)$ satisfies the iterated F -operator inequality, i.e.,

$$\tilde{Q}(x, u) \leq (F^M \tilde{Q})(x, u), \quad \forall x \in \mathcal{X}, u \in \mathcal{U},$$

and if $\hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$ allows for affine combinations of the basis functions, then \tilde{Q} is also an element of $\hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$.

Proof: Let $\beta \in \mathbb{R}$ denote the constant *downwards shift term* for notational convenience. Using the definition of the

F -operator we see that for any function $Q(x, u)$,

$$\begin{aligned} & (F(Q + \beta))(x, u) \\ & = l(x, u) + \gamma \min_{v \in \mathcal{U}} \mathbb{E}[Q(f(x, u, \xi), v) + \beta] \\ & = (FQ)(x, u) + \gamma \beta. \end{aligned}$$

where the equalities hold for all $x \in \mathcal{X}$, $u \in \mathcal{U}$. The first equality comes from the definition of the F -operator, and the second equality holds as β is an additive constant in the objective of the minimization.

Iterating the same argumentation M -times leads to

$$\begin{aligned} & (F^M(Q + \beta))(x, u) \\ & = (F^{M-1}(F(Q + \beta)))(x, u) \\ & = (F^{M-1}((FQ) + \gamma\beta))(x, u) \\ & = (F^{M-2}((F^2Q) + \gamma^2\beta))(x, u) \\ & = \dots \\ & = (F^M Q)(x, u) + \gamma^M \beta, \end{aligned} \quad (22)$$

where the equivalences hold point-wise for all $x \in \mathcal{X}$, $u \in \mathcal{U}$. Now we show that \tilde{Q} satisfies the iterated F -operator inequality,

$$\begin{aligned} & (F^M \tilde{Q})(x, u) \\ & = (F^M \hat{Q})(x, u) - \gamma^M \left(\frac{1 + \gamma^M}{1 - \gamma^M} \|Q^* - \hat{Q}\|_\infty \right) \\ & \geq \hat{Q}(x, u) - (1 + \gamma^M) \|Q^* - \hat{Q}\|_\infty \\ & \quad - \gamma^M \left(\frac{1 + \gamma^M}{1 - \gamma^M} \|Q^* - \hat{Q}\|_\infty \right) \\ & = \tilde{Q}(x, u), \end{aligned}$$

where the first equality comes from (22), the inequality is a direct application of Lemma 1.1 to the term $(F^M \hat{Q})$ and holds for all $x \in \mathcal{X}$, $u \in \mathcal{U}$, and the final equality follows from (21).

Finally, if $\hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$ allows for affine combinations of the basis functions, then $\hat{Q} \in \hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$ implies $\tilde{Q} \in \hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$ as the *downward shift term* is an additive constant. ■

Theorem 4.4 Given Q^* is the solution of (2), and \hat{Q}^* is the solution of (8) for a given choice $\hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$ and $c(\cdot, \cdot)$, then the following bound holds,

$$\|Q^* - \hat{Q}^*\|_{1, c(x, u)} \leq \frac{2}{1 - \gamma^M} \min_{\hat{Q} \in \hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})} \|Q^* - \hat{Q}\|_\infty$$

Proof: Given any approximate \mathcal{Q} -function from the basis, $\hat{Q}(x, u) \in \hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$, Lemma 1.2 allows us to construct,

$$\tilde{Q}(x, u) = \hat{Q}(x, u) - \frac{1 + \gamma^M}{1 - \gamma^M} \|Q^* - \hat{Q}\|_\infty \in \hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U}),$$

which is feasible for (8).

Working from the left hand side of the theorem,

$$\begin{aligned}
& \left\| Q^* - \hat{Q}^* \right\|_{1,c(x,u)} \\
& \leq \left\| Q^* - \hat{Q} \right\|_{1,c(x,u)} \\
& \leq \left\| Q^* - \hat{Q} \right\|_{\infty} \\
& \leq \left\| Q^* - \hat{Q} \right\|_{\infty} + \left\| \hat{Q} - \tilde{Q} \right\|_{\infty} \\
& = \left\| Q^* - \hat{Q} \right\|_{\infty} + \frac{1+\gamma^M}{1-\gamma^M} \left\| Q^* - \hat{Q} \right\|_{\infty} \\
& = \frac{2}{1-\gamma^M} \left\| Q^* - \hat{Q} \right\|_{\infty}
\end{aligned}$$

where the first inequality holds by Lemma 3.1 because \tilde{Q} is also feasible for (8), the second inequality by assuming w.l.o.g. that $c(x, u)$ is a probability distribution, the third inequality is an application of the triangle inequality, the first equality stems directly from the definition of \hat{Q} , and the final is an algebraic manipulation.

As this argumentation holds for any $\hat{Q} \in \hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$, the result follows. ■

C. Proof of Lyapunov-based fitting bound

The proof of Theorem 4.8 requires four auxiliary lemmas that are presented prior to the proof of Theorem 4.8. Lemma 1.3 bounds the difference after applying M iterations of the Bellman operator to 2 different Value functions. The bound is given by M iterations of the H_V operator introduced in Section IV-C and is used in Lemma 1.4 to give a bound on how much the M -iterated Bellman inequality is violated for any given Value function. This constraint violation bound is given in terms of a Lyapunov function and is used in Lemma 1.6 to prove that given any $\hat{V} \in \hat{\mathcal{F}}(\mathcal{X})$, it can be downshifted by a scalar multiple of a Lyapunov function to satisfy the M -iterated Bellman inequality. The Lyapunov function appearing in the downshift relates directly to the Lyapunov function and relevance weighting on the right-hand-side of the Theorem 4.8 bound.

Lemma 1.3: for any two functions $V_1, V_2 : \mathcal{X} \rightarrow \mathbb{R}$, and any integer $M \geq 1$,

$$|(\mathcal{T}^M V_1)(x) - (\mathcal{T}^M V_2)(x)| \leq \gamma^M (H_V^M(|V_1 - V_2|))(x),$$

for all $x \in \mathcal{X}$

Proof: The lemma will be proven by induction. For $M = 1$, we first show that the inequality hold without $|\cdot|$. Letting u_1^* denote the minimizer for $\mathcal{T}V_1$ and u_2^* for $\mathcal{T}V_2$,

$$\begin{aligned}
& (\mathcal{T}V_1)(x) - (\mathcal{T}V_2)(x) \\
& = (\mathcal{T}_u V_1)(x, u_1^*) - (\mathcal{T}_u V_2)(x, u_2^*) \\
& \leq (\mathcal{T}_u V_1)(x, u_2^*) - (\mathcal{T}_u V_2)(x, u_2^*) \\
& \leq \gamma \max_{u \in \mathcal{U}} ((\mathcal{T}_u V_1)(x, u) - (\mathcal{T}_u V_2)(x, u)) \\
& \leq \gamma \max_{u \in \mathcal{U}} |\mathbb{E}[V_1(f(x, u, \xi))] - \mathbb{E}[V_2(f(x, u, \xi))]|,
\end{aligned} \tag{23}$$

where the inequalities hold for all $x \in \mathcal{X}$. The first equality is the definition of \mathcal{T} in terms of \mathcal{T}_u , and the first inequality

holds by definition of u_1^* being the minimizer for $\mathcal{T}V_1$. The second inequality holds as the same u_2^* appears in both terms. The final inequality holds by definition of \mathcal{T}_u and $|\cdot|$.

An entirely analogous argument establishes that $(\mathcal{T}V_2)(x) - (\mathcal{T}V_1)(x)$ is bounded above by the same final term in (23). Hence the result for $M = 1$ follows as,

$$\begin{aligned}
& |(\mathcal{T}V_1)(x) - (\mathcal{T}V_2)(x)| \\
& \leq \gamma \max_{u \in \mathcal{U}} |\mathbb{E}[V_1(f(x, u, \xi))] - \mathbb{E}[V_2(f(x, u, \xi))]| \\
& \leq \gamma \max_{u \in \mathcal{U}} \mathbb{E}[|V_1(f(x, u, \xi)) - V_2(f(x, u, \xi))|] \\
& = \gamma (H_V(|V_1 - V_2|))(x),
\end{aligned} \tag{24}$$

where the inequalities hold for all $x \in \mathcal{X}$. The first inequality follows from (23). The second inequality uses [7, Lemma 1.7.2] to exchange the expectation and absolute value. The final equivalence is the definition of H_V as per Section IV-C.

Assume the statement holds true for some $k \in \mathbb{N}$, i.e.,

$$|(\mathcal{T}^k V_1)(x) - (\mathcal{T}^k V_2)(x)| \leq \gamma^k (H_V^k(|V_1 - V_2|))(x),$$

and show it therefore holds true for $k + 1$:

$$\begin{aligned}
& |(\mathcal{T}^{k+1} V_1)(x) - (\mathcal{T}^{k+1} V_2)(x)| \\
& = |(\mathcal{T}^k (\mathcal{T}V_1))(x) - (\mathcal{T}^k (\mathcal{T}V_2))(x)| \\
& \leq \gamma^k (H_V^k(|(\mathcal{T}V_1) - (\mathcal{T}V_2)|))(x) \\
& \leq \gamma^k (H_V^k(\gamma H_V(|V_1 - V_2|)))(x) \\
& = \gamma^{k+1} (H_V^{k+1}(|V_1 - V_2|))(x)
\end{aligned}$$

where the inequalities hold for all $x \in \mathcal{X}$. The first equivalence splits \mathcal{T}^{k+1} so that the induction assumption can be used to establish the first inequality. The second inequality uses (24) and the monotonicity property of H_V^k . The final equivalence follows by algebra.

By induction the claim holds for any integer $M \geq 1$. ■

Lemma 1.4: For any positive function $V^+ : \mathcal{X} \rightarrow \mathbb{R}_{++}$, any value function $V : \mathcal{X} \rightarrow \mathbb{R}$, and any integer $M \geq 1$,

$$V(x) - (\mathcal{T}^M V)(x) \leq (V^+(x) + \gamma^M (H_V^M V^+)(x)) \epsilon$$

for all $x \in \mathcal{X}$, where $\epsilon = \|V^* - V\|_{\infty, 1/V^+}$.

Proof: First we find a relation between V^+ , V , and V^* based on the weighted infinity norm.

$$\begin{aligned}
\epsilon V^+(x) & = \|V^* - V\|_{\infty, 1/V^+} V^+(x) \\
& \geq |V^*(x) - V(x)| (1/V^+(x)) V^+(x) \\
& = |V^*(x) - V(x)| \\
& \geq V(x) - V^*(x)
\end{aligned} \tag{25}$$

where the inequalities hold for all $x \in \mathcal{X}$. The first inequality comes from the definition of the weighted ∞ -norm. The first equality holds as V^+ is a strictly positive function, and the final inequality stems from the definition of $|\cdot|$.

Thus,

$$\begin{aligned}
& V(x) - (\mathcal{T}^M V)(x) \\
& \leq \epsilon V^+(x) + V^*(x) - (\mathcal{T}^M V)(x) \\
& \leq \epsilon V^+(x) + |(\mathcal{T}^M V^*)(x) - (\mathcal{T}^M V)(x)| \\
& \leq \epsilon V^+(x) + \gamma^M (H_V^M (|V^* - V|))(x) \\
& \leq \epsilon V^+(x) + \gamma^M (H_V^M (\epsilon V^+))(x) \\
& = \epsilon V^+(x) + \gamma^M \epsilon (H_V^M V^+)(x) \\
& = (V^+(x) + \gamma^M (H_V^M V^+)(x)) \epsilon
\end{aligned}$$

where the inequalities hold for all $x \in \mathcal{X}$. The first inequality is a consequence of (25). The second inequality uses the fact that $V^* = \mathcal{T}^M V^*$ and the definition of $|\cdot|$. The third inequality is a direct application of Lemma 1.3. The fourth inequality uses (25) and the monotonicity of operator H_V^M . The two equalities follow from simple algebra. ■

Lemma 1.5: Given any Lyapunov function V , as per definition 4.6, and its respective Lyapunov constant β_V , then,

$$\begin{aligned}
& \left(\frac{2}{1 - \beta_V^M} - 1 \right) (V(x) - \gamma^M (H_V^M V)(x)) \\
& \geq (V(x) + \gamma^M (H_V^M V)(x))
\end{aligned}$$

for all $x \in \mathcal{X}$.

Proof: By the definition of the Lyapunov function that $(HV)(x) \leq (\beta_V/\gamma) V(x)$ for all $x \in \mathcal{X}$, thus we get that,

$$\begin{aligned}
(H_V^M V)(x) &= (H_V^{M-1} (H_V V))(x) \\
&\leq (H_V^{M-1} ((\beta_V/\gamma)V))(x) \\
&= (\beta_V/\gamma) (H_V^{M-1} V)(x)
\end{aligned}$$

where the inequality holds for all $x \in \mathcal{X}$ by the monotone property of H^k for any $k \in \mathbb{N}$. Iterating the same argumentation M -times leads to,

$$(H_V^M V)(x) \leq (\beta_V/\gamma)^M V(x),$$

for all $x \in \mathcal{X}$. As V is strictly positive, this implies that,

$$\frac{2}{1 - \frac{\gamma^M (H_V^M V)(x)}{V(x)}} - 1 \leq \frac{2}{1 - \beta_V^M} - 1,$$

for all $x \in \mathcal{X}$. Manipulating the left-hand-side,

$$\left(\frac{2}{1 - \frac{\gamma^M (H_V^M V)(x)}{V(x)}} - 1 \right) = \frac{V(x) + \gamma^M (H_V^M V)(x)}{V(x) - \gamma^M (H_V^M V)(x)}$$

Hence the result follows. ■

Lemma 1.6: Let $\hat{V}^+(x) \in \hat{\mathcal{F}}(\mathcal{X})$ be a Lyapunov function, as per definition 4.6, and let $\hat{V} \in \hat{\mathcal{F}}(\mathcal{X})$ be an arbitrary element from the basis, and let \tilde{V} be defined as,

$$\tilde{V}(x) = \hat{V}(x) - \epsilon \left(\frac{2}{1 - \beta_{\hat{V}^+}^M} - 1 \right) \hat{V}^+(x) \quad (26)$$

where $\epsilon = \|V^* - \hat{V}\|_{\infty, 1/\hat{V}^+}$, then

$$\tilde{V}(x) \leq (\mathcal{T}^M \tilde{V})(x), \quad \forall x \in \mathcal{X}$$

i.e., it is feasible for the approximate iterated LP. Additionally \tilde{V} is an element of $\hat{\mathcal{F}}(\mathcal{X})$.

Proof: Starting from the right-hand-side of the iterated Bellman inequality,

$$\begin{aligned}
& (\mathcal{T}^M \tilde{V})(x) \\
&= (\mathcal{T}^M \hat{V})(x) - (\mathcal{T}^M \hat{V})(x) + (\mathcal{T}^M \tilde{V})(x) \\
&\geq (\mathcal{T}^M \hat{V})(x) - |(\mathcal{T}^M \hat{V})(x) - (\mathcal{T}^M \tilde{V})(x)| \\
&\geq (\mathcal{T}^M \hat{V})(x) - \gamma^M (H_V^M |\hat{V}(x) - \tilde{V}(x)|) \\
&= (\mathcal{T}^M \hat{V})(x) - \gamma^M \epsilon \left(\frac{2}{1 - \beta_{\hat{V}^+}^M} - 1 \right) (H^M \hat{V}^+)(x) \\
&\geq \hat{V}(x) - \epsilon (\hat{V}^+(x) + \gamma^M (H^M \hat{V}^+)(x)) \\
&\quad - \gamma^M \epsilon \left(\frac{2}{1 - \beta_{\hat{V}^+}^M} - 1 \right) (H^M \hat{V}^+)(x) \\
&= \tilde{V}(x) - \epsilon (\hat{V}^+(x) + \gamma^M (H^M \hat{V}^+)(x)) \\
&\quad + \epsilon \left(\frac{2}{1 - \beta_{\hat{V}^+}^M} - 1 \right) (\hat{V}^+(x) - \gamma^M (H^M \hat{V}^+)(x)) \\
&\geq \tilde{V}(x)
\end{aligned}$$

where the inequality holds for all $x \in \mathcal{X}$. The first equality is simple algebra and the first inequality is from the definition of $|\cdot|$. The second inequality is a direct application of Lemma 1.3. The second equality follows from the definition of \tilde{V} given in (26). The third inequality stems from applying Lemma 1.4 to the $(\mathcal{T}^M \hat{V})$ term. The last equality again uses the definition of \tilde{V} and the last inequality follows from Lemma 1.5.

By (26), \tilde{V} is a linear combination of \hat{V} and \hat{V}^+ . As \hat{V} and \hat{V}^+ are both elements of $\hat{\mathcal{F}}(\mathcal{X})$, so is \tilde{V} . ■

Theorem 4.8 Given V^* is the solution of (1), and \hat{V}^* is the solution of the approximate iterated LP. Then, for any function $\hat{V}^+(x) \in \hat{\mathcal{F}}(\mathcal{X})$ that is a Lyapunov Value function as per definition 4.6, the following bound holds,

$$\|V^* - \hat{V}^*\|_{1, c(x)} \leq \frac{2 \|\hat{V}^+\|_{1, c(x)}}{1 - \beta_{\hat{V}^+}^M} \min_{\hat{V} \in \hat{\mathcal{F}}} \|V^* - \hat{V}\|_{\infty, 1/\hat{V}^+}$$

Proof: Given any approximate Value function taken from the restricted function space, $\hat{V}(x) \in \hat{\mathcal{F}}(\mathcal{X})$, Lemma 1.6 allows us to construct $\tilde{V} \in \hat{\mathcal{F}}(\mathcal{X})$ as per (26) to be feasible for the approximate iterated LP, i.e., it satisfies $\tilde{V}(x) \leq (\mathcal{T}^M \tilde{V})(x)$ for all $x \in \mathcal{X}$. Working from the left-

hand-side of the bound,

$$\begin{aligned}
& \|V^* - \hat{V}^*\|_{1,c} \\
& \leq \|V^* - \tilde{V}\|_{1,c} \\
& = \int_{\mathcal{X}} \left(\frac{\hat{V}^+(x)}{\tilde{V}^+(x)} \right) |V^*(x) - \tilde{V}(x)| c(dx) \\
& \leq \left(\int_{\mathcal{X}} \hat{V}^+(x) c(dx) \right) \max_{z \in \mathcal{X}} \frac{|V^*(z) - \tilde{V}(z)|}{\hat{V}^+(z)} \\
& = \left(\|\hat{V}^+\|_{1,c(x)} \right) \|V^* - \tilde{V}\|_{\infty, 1/\hat{V}^+} \\
& \leq \|\hat{V}^+\|_{1,c(x)} \left(\|V^* - \hat{V}\|_{\infty, 1/\hat{V}^+} + \|\hat{V} - \tilde{V}\|_{\infty, 1/\hat{V}^+} \right) \\
& = \|\hat{V}^+\|_{1,c(x)} \left(\frac{2}{1 - \beta_{\hat{V}^+}^M} \right) \|V^* - \hat{V}\|_{\infty, 1/\hat{V}^+}
\end{aligned}$$

where the inequalities hold for all $x \in \mathcal{X}$. The first inequality follows from Lemma 3.1 and Lemma 1.6. The first equality is the definition of the weighted 1-norm and holds as \hat{V}^+ is strictly positive. The second inequality holds because the objective of the maximization is non-negative for all $z \in \mathcal{X}$. The second equality is the definition of the weighted 1-norm and weighted ∞ -norm. The final inequality follows by the triangle inequality. The final equality stems from using (26) by taking the weighted ∞ -norm of $(\hat{V} - \tilde{V})$ and then some simple algebra.

As the inequality established holds for any $\hat{V}(x) \in \hat{\mathcal{F}}(\mathcal{X})$, it also holds when the right-hand-side is minimized over all $\hat{V}(x) \in \hat{\mathcal{F}}(\mathcal{X})$. Hence the result follows. ■

D. Proofs of equivalent \mathcal{Q} -function formulation

Lemma 5.1: If sets $\hat{\mathcal{F}}(\mathcal{X})$ and $\hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$ satisfy that for all $\hat{V} \in \hat{\mathcal{F}}(\mathcal{X})$ there exists a $\hat{Q} \in \hat{\mathcal{F}}(\mathcal{X} \times \mathcal{U})$ such that

$$\hat{Q}(x, u) = \mathcal{T}_u \hat{V}(x, u), \quad \forall x \in \mathcal{X}, u \in \mathcal{U},$$

then the approximate LP (13) is equivalent to (14), in the sense that both problems have the same optimal value and there is a mapping between feasible and optimal solutions in both problems.

Proof: We shall show that any feasible solution of (13) corresponds a feasible solution of (14) with the same objective value, and vice versa. Note that for the proof superscript $(\cdot)'$ indicates a decision variable of problem (14).

Suppose that $\{\hat{Q}_j\}_{j=0}^M, \{\hat{V}_j\}_{j=0}^{M-1}$ is a feasible solution of (13), and take the following decision variables for (14),

$$\hat{Q}'_0 = \hat{Q}_0, \quad \hat{V}'_j = \hat{V}_j, \quad j = 0, \dots, M-1.$$

Feasibility for the constraints of (14) are now checked.

$$\hat{Q}'_0(x, u) = \hat{Q}_0(x, u) \leq \mathcal{T}_u \hat{V}_0(x, u) = \mathcal{T}_u \hat{V}'_0(x, u),$$

for all $x \in \mathcal{X}$ and $u \in \mathcal{U}$, thus (14a) is satisfied. For $j = 1, \dots, M-1$ we have that,

$$\hat{V}'_{j-1}(x) = \hat{V}_{j-1}(x) \leq \hat{Q}'_j(x, u) \leq \mathcal{T}_u \hat{V}_j(x, u) = \mathcal{T}_u \hat{V}'_j(x, u),$$

for all $x \in \mathcal{X}$ and $u \in \mathcal{U}$, thus (14b) are satisfied. Finally,

$$\hat{V}'_{M-1}(x) = \hat{V}_{M-1}(x) \leq \hat{Q}_M(x, u) = \hat{Q}_0(x, u) = \hat{Q}'_0(x, u),$$

for all $x \in \mathcal{X}$ and $u \in \mathcal{U}$. Thus (14c) is also satisfied, and the considered decision variables are feasible for problem (14). As $\hat{Q}'_0 = \hat{Q}_0$, the objective values are equal.

This completes the equivalence in one direction. Now suppose that $\hat{Q}'_0, \{\hat{V}'_j\}_{j=0}^{M-1}$ is a feasible solution of (14), and take the following decision variables for (13) to be defined as,

$$\begin{aligned}
\hat{Q}_0 &= \hat{Q}_M = \hat{Q}'_0, \\
\hat{V}_j &= \hat{V}'_j, \quad j = 0, \dots, M-1 \\
\hat{Q}_j &= \mathcal{T}_u \hat{V}'_j, \quad j = 1, \dots, M-1,
\end{aligned}$$

where the choices of \hat{Q}_j are valid by the assumption of the lemma.

Feasibility for the constraints of (13) is now checked.

$$\hat{Q}_0(x, u) = \hat{Q}'_0(x, u) \leq \mathcal{T}_u \hat{V}'_0(x, u) = \mathcal{T}_u \hat{V}_0(x, u),$$

for all $x \in \mathcal{X}$ and $u \in \mathcal{U}$, and for $j = 1, \dots, M-1$ we have that,

$$\hat{Q}_j(x, u) = \mathcal{T}_u \hat{V}'_j(x, u) \leq \mathcal{T}_u \hat{V}_j(x, u),$$

for all $x \in \mathcal{X}$ and $u \in \mathcal{U}$, thus (13a) are satisfied. For $j = 0, \dots, M-2$ we also have that,

$$\hat{V}_j(x) = \hat{V}'_j(x) \leq \mathcal{T}_u \hat{V}'_{j+1}(x, u) = \hat{Q}_{j+1}(x, u),$$

for all $x \in \mathcal{X}$ and $u \in \mathcal{U}$, and for $j = M-1$ we have that,

$$\hat{V}_{M-1}(x) = \hat{V}'_{M-1}(x) \leq \hat{Q}'_0(x, u) = \hat{Q}_M(x, u),$$

for all $x \in \mathcal{X}$ and $u \in \mathcal{U}$, thus (13b) are satisfied. Constraint (13c) is satisfied by construction and thus the considered decision variables are feasible for problem (13). As $\hat{Q}_0 = \hat{Q}'_0$, the objective values are equal.

This completes the proof that problems (13) and (14) are equivalent in the sense proposed. ■

REFERENCES

- [1] Richard E Bellman. On the theory of dynamic programming. *Proceedings of the National Academy of Sciences of the United States of America*, 38(8):716–719, 1952.
- [2] A. Ben-Tal, A. Goryashko, E. Guslitzer, and A. Nemirovski. Adjustable robust solutions of uncertain linear programs. *Mathematical Programming*, 99(2):351–376, 2004.
- [3] Dimitri P Bertsekas. *Dynamic programming and optimal control*, volume 1. Athena Scientific Belmont, MA, 2005.
- [4] Paul Beuchat, Angelos Georgiou, and John Lygeros. Alleviating tuning sensitivity in approximate dynamic programming. In *Submitted to ECC 2016, but decision is not yet received*, 2016.
- [5] Francesco Borrelli, Alberto Bemporad, and Manfred Morari. *Predictive Control for linear and hybrid systems*. March 2014.
- [6] Eduardo F Camacho and Carlos Bordons Alba. *Model Predictive Control*. Springer Science & Business Media, 2007.
- [7] Ole Christensen. *Functions, spaces, and expansions: mathematical tools in physics and engineering*. Springer Science & Business Media, 2010.

- [8] Randy Cogill, Michael Rotkowitz, Benjamin Van Roy, and Sanjay Lall. An approximate dynamic programming approach to decentralized control of stochastic systems. In *Control of Uncertain Systems: Modelling, Approximation, and Design*, pages 243–256. Springer, 2006.
- [9] Daniela P De Farias and Benjamin Van Roy. The linear programming approach to approximate dynamic programming. *Operations Research*, 51(6):850–865, November-December 2003.
- [10] Daniela P De Farias and Benjamin Van Roy. On constraint sampling in the linear programming approach to approximate dynamic programming. *INFORMS - Mathematics of Operations Research*, 29(3):462–478, August 2004.
- [11] Angelos Georgioulou, Wolfram Wiesemann, and Daniel Kuhn. Generalized decision rule approximations for stochastic programming via liftings. *Mathematical Programming*, pages 1–38, 2010.
- [12] Onésimo Hernández-Lerma and Jean B Lasserre. *Discrete-time Markov control processes: basic optimality criteria*, volume 30. Springer Science & Business Media, 2012.
- [13] Nikolaos Kariotoglou, Sean Summers, Tyler Summers, Maryam Kamgarpour, and John Lygeros. Approximate dynamic programming for stochastic reachability. In *Control Conference (ECC), 2013 European*, pages 584–589. IEEE, 2013.
- [14] Arezou Keshavarz and Stephen Boyd. Quadratic approximate dynamic programming for input-affine systems. *International Journal of Robust and Nonlinear Control*, 24(3):432–449, July 2012.
- [15] K Krishnamoorthy, Meir Pachter, Swaroop Darbha, and Phil Chandler. Approximate dynamic programming with state aggregation applied to uav perimeter patrol. *International Journal of Robust and Nonlinear Control*, 21(12):1396–1409, 2011.
- [16] Warren B Powell. What you should know about approximate dynamic programming. *Naval Research Logistics (NRL)*, 56(3):239–249, February 2009.
- [17] Warren B Powell. Clearing the jungle of stochastic optimization. *Informa Tutorials*, 2014.
- [18] Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2005.
- [19] James B Rawlings and David Q Mayne. *Model Predictive Control: Theory and Design*. Nob Hill Publishing, 2009.
- [20] Paul J Schweitzer and Abraham Seidmann. Generalized polynomial approximations in markovian decision processes. *Journal of Mathematical Analysis and Applications*, 110:568–582, 1985.
- [21] Alexander Shapiro, Darinka Dentcheva, and Andrzej Ruszczyński. *Lectures on stochastic programming: modeling and theory*, volume 16. SIAM, 2014.
- [22] Tyler Summers, Konstantin Kunz, Nikolaos Kariotoglou, Maryam Kamgarpour, Sean Summers, and John Lygeros. Approximate dynamic programming via sum of squares programming. In *Control Conference (ECC), 2013 European*, pages 191–197. IEEE, 2013.
- [23] Tobias Sutter, Peyman Mohajerin Esfahani, and John Lygeros. Approximation of constrained average cost markov control processes. In *Decision and Control (CDC), 2014 IEEE 53rd Annual Conference on*, pages 6597–6602. IEEE, 2014.
- [24] Yang Wang, Brendan O’Donoghue, and Stephen Boyd. Approximate dynamic programming via iterated bellman inequalities. *International Journal of Robust and Nonlinear Control*, 2014.
- [25] Christopher Watkins. *Learning from Delayed Rewards*. PhD thesis, King’s College, May 1989.